

QuickTime™ and a
TIFF (LZW) decompressor
are needed to see this picture.

Strategies for Manual Annotation of the Anopheles Genome

Kathryn Campbell

VectorBase



Anopheles gambiae

Manual Gene Annotation

Goal: Significantly improve the quality of gene annotations to reveal strategies for targeting *mosquito:parasite:human* interactions

- 278 Mb genome sequenced in 2002
- 14,707 Ensembl automated gene predictions
- 140,000 ESTs

The primary vector of malaria in Sub-Saharan Africa, the mosquito *A. gambiae* is a major contributor to the more than half-a-billion cases of malaria each year



Using *Drosophila melanogaster* as a paradigm, manual annotation of the *Anopheles* genome is expected to greatly improve the quality of annotations due to limitations of gene prediction algorithms.



Strategy: Improve Annotations that represent computational challenges for automated gene builds

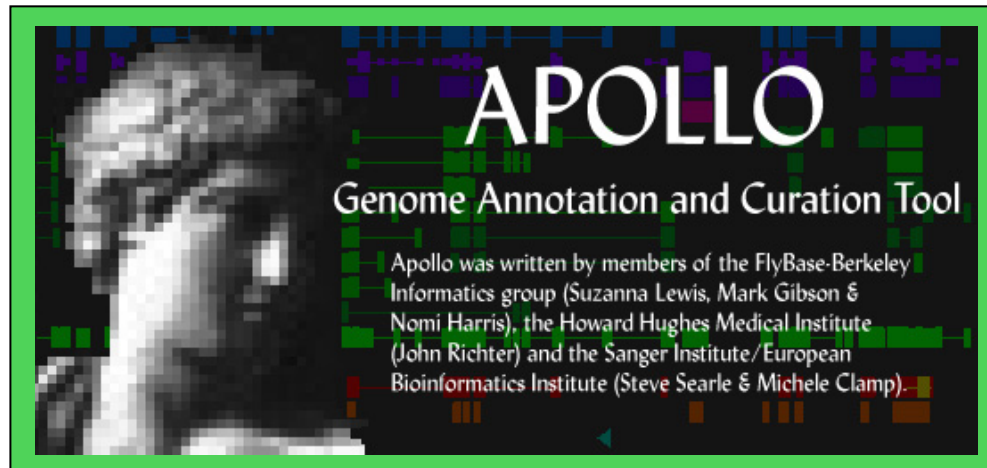
- **Alternative transcripts**
- **Overlapping genes**
- **Tandem homologous genes**
- **Dicistronic genes**
- **Nested genes**
- **Small ORFs**
- **Genes with repetitive sequence**

A first pass manual annotation will target regions and gene families of interest to the scientific community



Prioritizing Regions to be Annotated

- Genes involved in immune response, host preference, and insecticide resistance
- Multi-gene family members
- Tandem related genes
- 2 alternative assemblies of the same region due to divergent haplotypes in the PEST strain sequenced



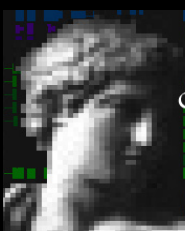
- Apollo is the graphical interface that will be used for editing annotations on the *Anopheles* genomic sequence
- Apollo was used by FlyBase to manually annotate the entire 117Mb *Drosophila melanogaster* Euchromatin (10 curators/7 months)
- Annotation editor for GMOD (Generic Model Organism Database)
- Developed by FlyBase-BDGP and Ensembl

Data Types used for Manual Annotation

- Ensembl - Automated gene build based on synthesis of known proteins and cDNA/EST sequences
- SNAP- Gene finder trained on Anopheles
- TGE_gw - Anopheles/Drosophila protein aligned to the genome
- Genewise - Uniprot protein sequences aligned to the genome
- Genomewise - Ensembl EST transcript predictions
- SWALL - BLASTX similarity of SWISS-PROT/TrEMBL dataset
- Drosophila-Peptides - blast hit by Drosophila proteins.
- RNA_BEST - Anopheles EST/mRNA
- Protein features - Pfam domains, Prints family fingerprints, transmembrane regions, signal peptides
- RepeatMasker, TRF, SEG- identifies repeats/low complexity features



Types			
<input type="checkbox"/> Sort	ensembl	Show	<input checked="" type="checkbox"/>
<input checked="" type="checkbox"/> Label		Expand	<input checked="" type="checkbox"/>
<input type="checkbox"/> Sort	Snap	Show	<input checked="" type="checkbox"/>
<input type="checkbox"/> Label		Expand	<input type="checkbox"/>
<input type="checkbox"/> Sort	Swall	Show	<input checked="" type="checkbox"/>
<input type="checkbox"/> Label		Expand	<input type="checkbox"/>
<input type="checkbox"/> Sort	drosophila-peptides	Show	<input checked="" type="checkbox"/>
<input type="checkbox"/> Label		Expand	<input type="checkbox"/>
<input type="checkbox"/> Sort	TGE_gw	Show	<input checked="" type="checkbox"/>
<input type="checkbox"/> Label		Expand	<input type="checkbox"/>
<input type="checkbox"/> Sort	similarity_genewise	Show	<input checked="" type="checkbox"/>
<input type="checkbox"/> Label		Expand	<input checked="" type="checkbox"/>
<input type="checkbox"/> Sort	CpG	Show	<input checked="" type="checkbox"/>
<input type="checkbox"/> Label		Expand	<input type="checkbox"/>
<input type="checkbox"/> Sort	RepeatMask	Show	<input type="checkbox"/>
<input type="checkbox"/> Label		Expand	<input type="checkbox"/>
<input type="checkbox"/> Sort	TRF	Show	<input type="checkbox"/>
<input type="checkbox"/> Label		Expand	<input type="checkbox"/>
<input type="checkbox"/> Sort	Dust	Show	<input type="checkbox"/>
<input type="checkbox"/> Label		Expand	<input type="checkbox"/>
<input type="checkbox"/> Sort	DnaAligns(Rna)	Show	<input checked="" type="checkbox"/>
<input type="checkbox"/> Label		Expand	<input checked="" type="checkbox"/>
<input type="checkbox"/> Sort	DnaAligns(RnaBest)	Show	<input checked="" type="checkbox"/>
<input type="checkbox"/> Label		Expand	<input checked="" type="checkbox"/>



APOLLO

Genome Annotation and Curation Tool

Apollo was written by members of the FlyBase-Berkeley Informatics group (Suzanna Lewis, Mark Gibson & Nomi Harris), the Howard Hughes Medical Institute (John Richter) and the Sanger Institute/European Bioinformatics Institute (Steve Searle & Michele Clamp).

VectorBase curators will utilize the APOLLO graphic interface to view and synthesize data from multiple regions of the Anopheles genome

File

Edit

View

Tiers

Analysis

Bookmarks

Annotation

Window

Links

Help

Type	Name	Range	Score
Ense...	Ense...	15501...	100.0

Ensembl_2: ensembl:ensembl...

Ensembl_2: ENSANGT00000021869

Name	Genomic ...	Phase
Ensembl_2_1550189...		0
Ensembl_2_1550161...		0

Chromosome

<

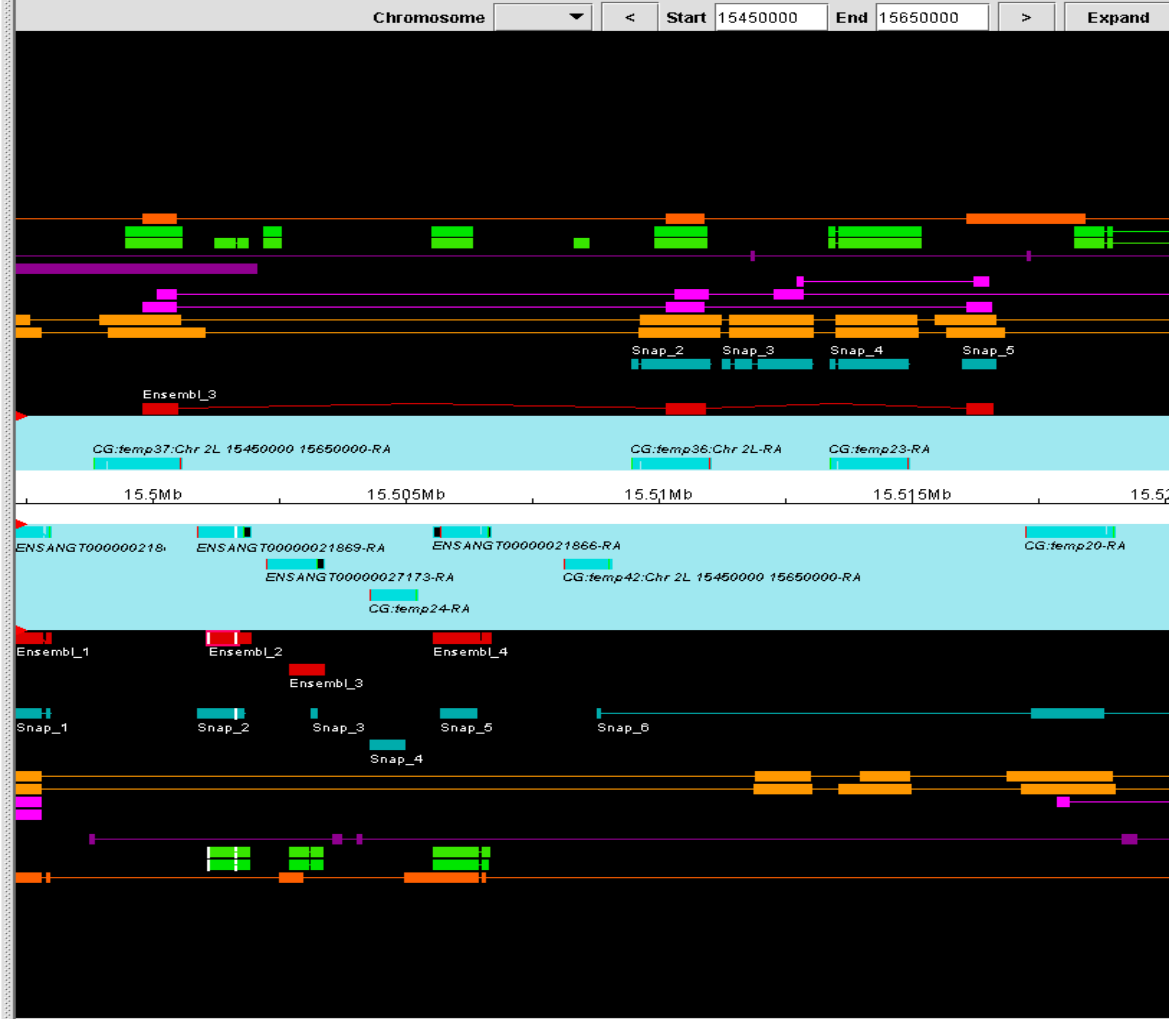
Start 15450000

End 15650000

>

Expand

Load



☒ Sort
 ☒ Label

ensembl

Show

Expand

☐ Sort
 ☒ Label

Snap

Show

Expand

☐ Sort
 ☐ Label

Swall

Show

Expand

☐ Sort
 ☐ Label

drosophila-peptides

Show

Expand

☐ Sort
 ☐ Label

TGE_gw

Show

Expand

☒ Sort
 ☒ Label

similarity_genewise

Show

Expand

☐ Sort
 ☐ Label

CpG

Show

Expand

☐ Sort
 ☐ Label

RepeatMask

Show

Expand

☐ Sort
 ☐ Label

TRF

Show

Expand

☐ Sort
 ☐ Label

Dust

Show

Expand

☐ Sort
 ☐ Label

DnaAligns(Rna)

Show

Expand

☐ Sort
 ☐ Label

DnaAligns(RnaBest)

Show

Expand

☐ Sort
 ☐ Label

ProteinAnnotation

Show

Expand

Position

15498654

Feature

Action

Zoom

x10

x2

x5

x.1

Reset

Zoom factor = 6.4000

A curator is able to synthesize complex data sets resulting in significantly improved gene annotations

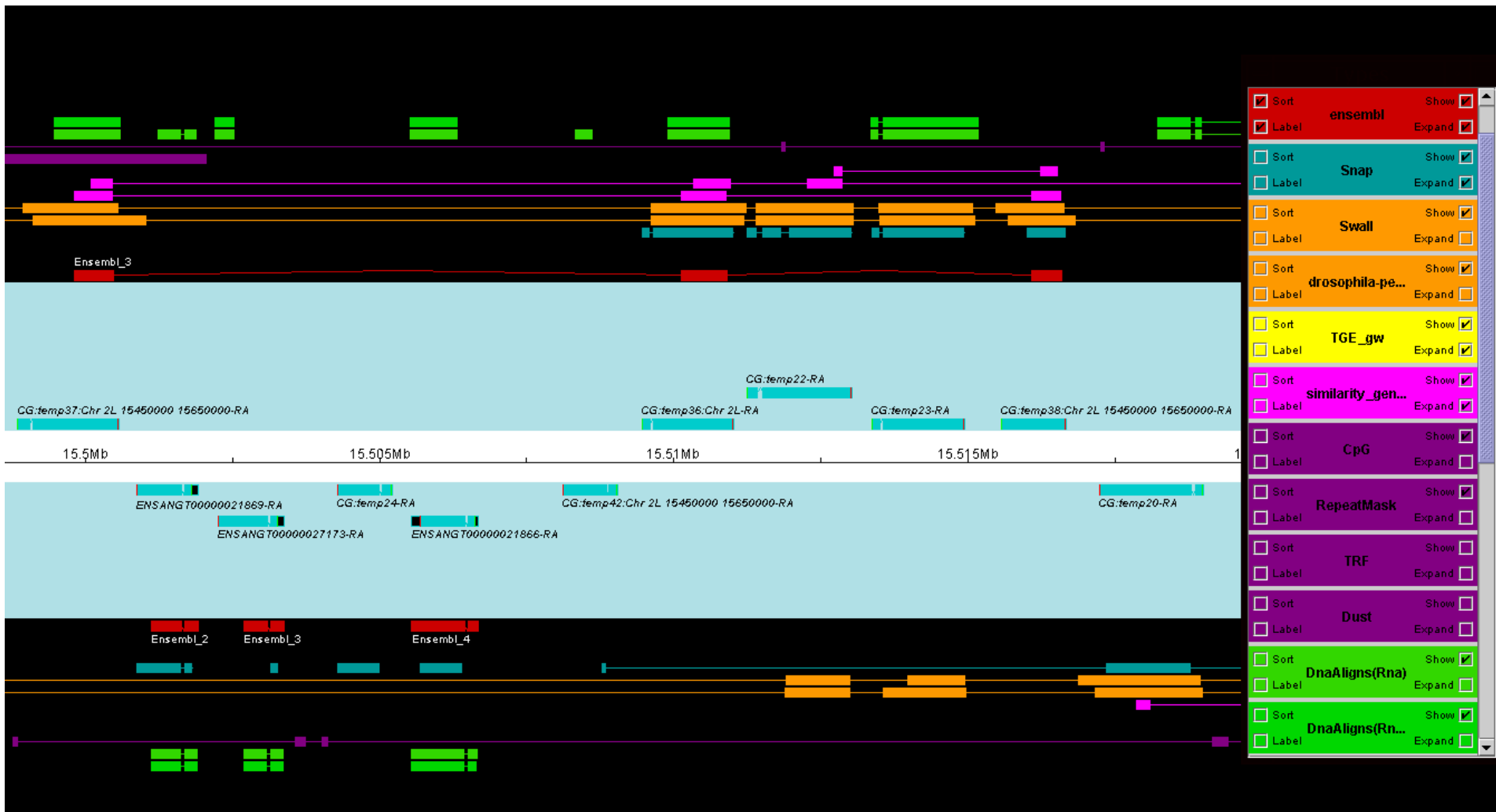
Types

<input checked="" type="checkbox"/> Sort	ensembl	Show <input checked="" type="checkbox"/>
<input checked="" type="checkbox"/> Label		Expand <input checked="" type="checkbox"/>
<input type="checkbox"/> Sort	Snap	Show <input checked="" type="checkbox"/>
<input type="checkbox"/> Label		Expand <input type="checkbox"/>
<input type="checkbox"/> Sort	Swal	Show <input checked="" type="checkbox"/>
<input type="checkbox"/> Label		Expand <input type="checkbox"/>
<input type="checkbox"/> Sort	drosophila-peptides	Show <input checked="" type="checkbox"/>
<input type="checkbox"/> Label		Expand <input type="checkbox"/>
<input type="checkbox"/> Sort	TGE_gw	Show <input checked="" type="checkbox"/>
<input type="checkbox"/> Label		Expand <input type="checkbox"/>
<input type="checkbox"/> Sort	similarity_genewise	Show <input checked="" type="checkbox"/>
<input type="checkbox"/> Label		Expand <input checked="" type="checkbox"/>
<input type="checkbox"/> Sort	CpG	Show <input checked="" type="checkbox"/>
<input type="checkbox"/> Label		Expand <input type="checkbox"/>
<input type="checkbox"/> Sort	RepeatMask	Show <input type="checkbox"/>
<input type="checkbox"/> Label		Expand <input type="checkbox"/>
<input type="checkbox"/> Sort	TRF	Show <input type="checkbox"/>
<input type="checkbox"/> Label		Expand <input type="checkbox"/>
<input type="checkbox"/> Sort	Dust	Show <input type="checkbox"/>
<input type="checkbox"/> Label		Expand <input type="checkbox"/>
<input type="checkbox"/> Sort	DnaAligns(Rna)	Show <input checked="" type="checkbox"/>
<input type="checkbox"/> Label		Expand <input checked="" type="checkbox"/>
<input type="checkbox"/> Sort	DnaAligns(RnaBest)	Show <input checked="" type="checkbox"/>
<input type="checkbox"/> Label		Expand <input checked="" type="checkbox"/>

Initial analysis: Eight 200-250 Kb regions



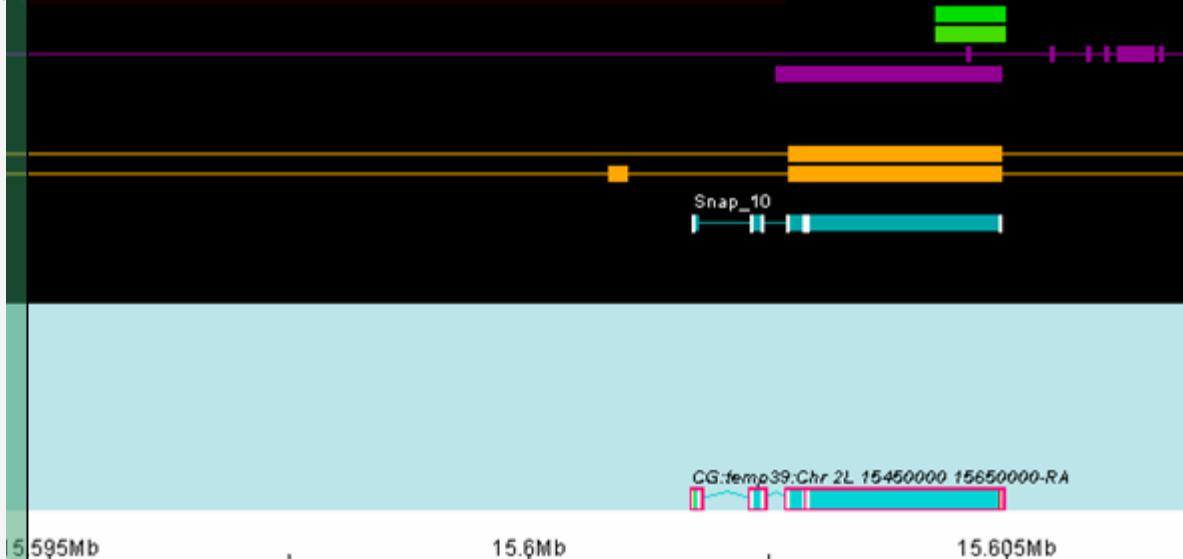
A Major Challenge for Automated Gene Annotation: Predicting Tandem Related Genes



Low Complexity sequences may require manual annotation

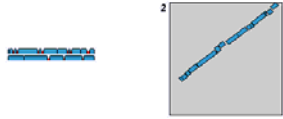
[illegible]

XCG:temp39:Chr 2L 15450000 15650000-RA (amino acid sequence)
 MLVDRDQLQERFEFRVSRTRRRCCDEGSSNRHTAGRLCLGSKFASRAARL
 YHHCILHIGRIYDFLCANYTAFDQKTFICHFASEVDCXNSPKYVFRNEPLYK
 ATTTTTVKPPPTTSTTAAPTTVPSRPKPLRKPVRRRRPQVDYVYDYEDYE
 EDYVEERNRKRKNRPNRRPLVDYDEDERFORVRYRERERERVEEDDOY
 VEDRRPYRMKTRPRNGQDRRPLYDDORRYDDORRRDRDRRPAVERPMRR
 PMPLODEVEERKPVAAVNERRRPSRPTVDSRRQAVGEDRRSPAQDAESRRG
 SVVEERKPAAKRPAYDORRYQDEVDDEEDDOYDARGSYGGGGGRRRAADKG
 RQQNDVMTVKPSGSSIIYDRPAAPRINRPVPLNEKSYAYYNPDAAGGKQ
 GAGKSTTAAPEGETYDEEYVPPARRDQKKEADGKPPRSAGSDVVYVYD
 VIVAKEDLPLRTAPGSVRKPEQSERVNGRONS ELAKGSGRGGIKHQPQSA
 PAPSKPNRNVLPERVEEEDAEYDVPAPRPSVIARQQAALYNFNSKSKSQ
 QREPPTTESASTARAPVRSTKRPFLPSRGGNPYLARGLQPVGVAKQFAGG
 KSSPEKPAFPRIDIDGATTPQAVSGKRQQQQAPPAAPSPAQPPQQQQQQ
 QQEVNVNVTLDQLYDDEYDVTLNDAINPTLKLPTRSAPQNFRNRVYHQ
 DEVFVPEFAPAYPSEFRRTAIRPPAPLLYHGSVAIEAGSANTSQRRGAH
 VFGTYEY



Finding novel genes requires manual annotation

Sequence 1 lcl|_CG:temp32-Chr 2L 30990000 31240000-RB (amino acid sequence): 443 residues Length 443
Sequence 2 lcl|_CG:temp1.CG12090-RA (amino acid sequence): 592 residues Length 592

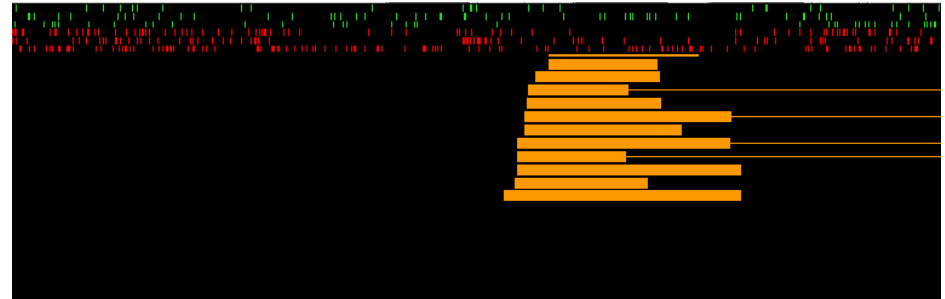


NOTE: The statistics (bitscore and expect value) is calculated based on the size of nr database

Score = 47.8 bits (112), Expect = 0.001
Identities = 87/422 (20%), Positives = 163/422 (38%), Gaps = 68/422 (16%)



Query: 46 LFPFPLR-TLGNKINALLSNYP----YATLLSSVME--GVDAVILTARESFHFEVOL 98
+PPP R G+++ + P Y + ME G+ ++ + L
Sbjct: 191 VTFPANKRNFGRFRLFLVQGVDPNTFWYKRRKAMELAGLGGILINQMSGLAVTMDL 250
Query: 99 INHAINRQNALHMEILROBLENGTIDHMFRAHVSRLNVDAAYMEALVAVPKS 158
+N + L+ M + + G +++ +S +VD + + + +P
Sbjct: 251 FFEVNGSSLLN-MAALTDLVKGKVELSPHLVDTLQNTSDVSYTPTVAPRCFMPLD 309
Query: 159 DQNLINIMLQPTTIVNTIVLAYLL-----VGMIKLSPFPK-----RRCNRP 202
++++ + PF++ +W +L LL VR++I P+ R NR
Sbjct: 310 NEIRSLVFLPFLSTHMLCLLPVLLVHFVTVRLIPDGHFMAILVPGAGQVRYGNRK 369
Query: 203 TLKKTLCRMSFGSUSLTTVVVELVSPFLIEAYLAKITEFLLYCRFRSDPQTLDEFFRS 262
+++ T PG F+L + Y K+T L R +L+R F
Sbjct: 370 FVRRPFTFLILPG-----IFLQQTITSLTSLVILIRPDFWSELEELFLL 416
Query: 263 TIPVLV-PETMDLVEALOPTVAANPHAKLIRPDYAKPAASCCARHTLPRAEVVR-- 319
+DV P + +V++G + +++++ EI P Y +
Sbjct: 417 FTRILVLPDVVAIVDSLGHAEQFSTKFSCTDAENFSQK-----RISMHPETIIPSTI 470
Query: 320 NKKYFDAT---LGRKQLYLPEQLTIIPMSYLWGRFAKNSFELFVHESCLIGTV 376
++PD L +K+ Y P Y + + K+ FLL V ++OL ++
Sbjct: 471 PWRFFDMQGRFLKKRFTYFKICHGSPFYQQLRVDSHLKDALHRLFLHQQAGLNDLW 530
Query: 377 -TAMRKDMMHMERHFFTKGMLTAD-----LLPLFLVAGWCCSFAAFLE 421
T+RK R + K + TL+ L+P F L G S AFL
Sbjct: 531 DTCYRK---ARRWYLEDFTLAELKEELRFLPLALNLLVPAPSLFLCOMLSGIAFLV 586
Query: 422 EV 423
E+
Sbjct: 587 EI 588



CG:temp32-Chr 2L 30990000 31240000-RB

ENSANGP00000029595
ENSANGP00000025936
CG_temp32_Ch2L309900003124000
ENSANGP00000028419
ENSANGP00000026593

VATRDMVVLDDHA--PKLHLTSRDYFNLVVPKPVKLNLDAMVQPFNGTVW 263
VATRDMVGLDHT--PKLHLTSRDYFNLVVPKPVKLNLDAMVQPFNGTVW 342
HMPRAHVSRSYLVNDVAAAYEWEALVVLVVPKSDQLNLINIMLQPTTIEVW 176
MMTRRHVNVGIL--PVVYIPDITYYCLVAPRTTQIDLQSLRPFSGTVW 228
DLTFATSTFRAHYMHDVSLKERGGYCVLCPFHTEPDLRHLKLPFSFGIW 214
.: * : : : : : *

ENSANGP00000029595
ENSANGP00000025936
CG_temp32_Ch2L309900003124000
ENSANGP00000028419
ENSANGP00000026593

LMIGALLGVRFVAGYQLQDLFGMLDRSGTLRKWLRLSQALHFP--CPQWVQ 311
LMIGVLLGVRFVAGYQLQDLFGMLDRSGTKVRKWRSLQAPHCP--CTQWVQ 390
TIVLAALLVQMI---KLFSFFKRRCNRFLLKTLTCRWSFGSFGSLTT 222
WFIVVCTVL-----ISAFDELKQHTRLGRLAQQLFARQPIASFYR 269
AVLGALLVGCRLLG---HLFPALFERNLEQIFFTAGASHRQPFPTTRIVS 261
.: : :

ENSANGP00000029595
ENSANGP00000025936
CG_temp32_Ch2L309900003124000
ENSANGP00000028419
ENSANGP00000026593

LAADVLTFLLEIAYLAQVTSLLLTLRFIEGPRDLNEFIASNIRIVEPYES 361
LAADVLTFLLEIAYLAQVTSLLLTLRFIEGPRDLNEFIASNMIRIVEPYES 440
VGVELVSFLLIEAYLAKITEFLLYCRFRSDPQTLDEFFRSTIPVLVPEYM 272
ICLAVISFVLIESYLATVTSFFLAIRFVDPDAKLEEFFATGIPIRLPEGM 319
FSAAVLIFLSEAYNAKIVSLMSDSKYFDRPESVRELIESDLKVAIPGVR 311
.: * * * * : : : : : * :

ENSANGP00000029595
ENSANGP00000025936
CG_temp32_Ch2L309900003124000
ENSANGP00000028419
ENSANGP00000026593

TLSLVQVETGQORALLKTRFVKRTAEELARPNOQDAFVELSRVAFQSYGT 411
TLSLVQVETGQORALLKTRFVKRTAEELARPNOQDAFVELSRVAFQSYGT 490
DPLVEALGPTVAANFHAHLIR-PDEYAKRAASCCARHTLPRAEVVRMG 321
VRFLRNLEPLKDRIMARGVG--ASVCSPISTQCAHLESFAMASYQISEN 367
ASLLAESLPGKLVNRRR-----AEALYRERGPLLNFNEYCTVMYCYLAYLQ 356
.: . :

ENSANGP00000029595
ENSANGP00000025936
CG_temp32_Ch2L309900003124000
ENSANGP00000028419
ENSANGP00000026593

EPFDPVTGRNRFYLLDEPLADVR-FQYSFAKETAFMGVAQDYMLRCEENG 460
EPFDPVTGRNRFYLLDEPLADVR-FQYSFAKETAFMGVAQDYMLRCEENG 539
KYFDA TLGRKQLYILPEQLTIIP-MSYLVGRNFAFNKSFELFLLSVHESG 370
VGIDPVSGRKRSYIVPEMVTIRYSYLSYAFARGSPDLTVAMYLWRMYEAG 417
TSVGRNLHGQQYVLPDMVSBQL-RTLQATHSPFYDTFAEYERYFQSG 405
.. * : : : : : : : *

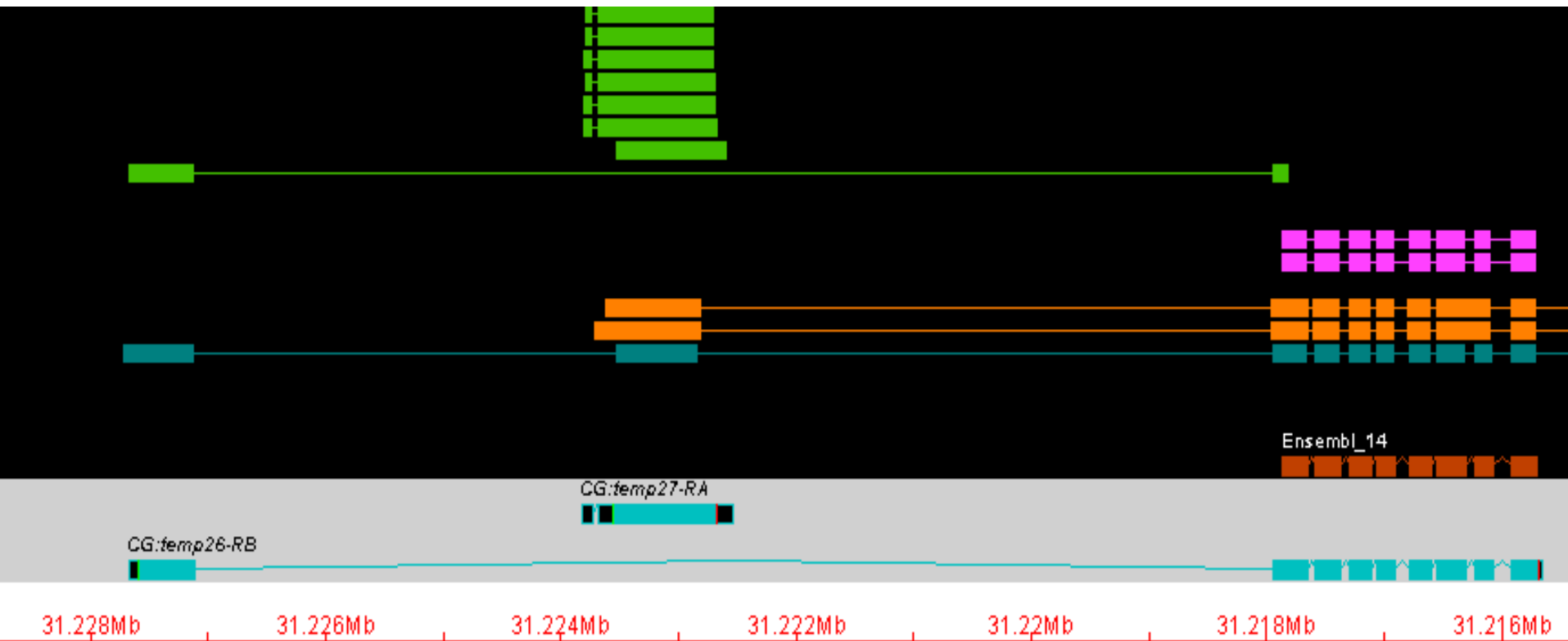
0.994Mb 30.996Mb 30.998Mb

CG:temp2-RA

Ensembl_1

Snap_14

Nested genes pose a challenge for gene prediction algorithms



Putative Anopheles specific gene identified by pattern recognition and visual search for ORFs

Sequences producing significant alignments:			Score (bits)	E Value
gi 55238289 gb EAA11136.3 	ENSANGP000000021866	[Anopheles ga...	229	7e-59
gi 55238287 gb EAA11135.2 	ENSANGP000000021869	[Anopheles ga...	177	2e-43
gi 55238288 gb EAL39836.1 	ENSANGP000000027562	[Anopheles ga...	92	2e-17
gi 55239049 gb EAL40035.1 	ENSANGP000000025919	[Anopheles ga...	75	2e-12
gi 55238124 gb EAL39790.1 	ENSANGP000000025599	[Anopheles ga...	70	7e-11
gi 55240053 gb EAL40329.1 	ENSANGP000000029231	[Anopheles ga...	60	6e-08
gi 55244958 gb EAL41642.1 	ENSANGP000000026714	[Anopheles ga...	56	1e-06
gi 55241665 gb EAL40761.1 	ENSANGP000000028587	[Anopheles ga...	48	4e-04
gi 55244020 gb EAL41360.1 	ENSANGP000000029329	[Anopheles ga...	38	0.37
gi 49646778 emb CAG83163.1 	unnamed protein product [Yarrow...		36	1.5
gi 297855 emb CAA42211.1 	fatty-acid synthase [Yarrowia lip...		35	1.6
gi 56468200 gb EAL46075.1 	phosphatidylinositol 3-kinase, p...		34	4.9
gi 56459878 ref YP_155159.1 	Thiamine biosynthesis protein ...		34	4.9
gi 34067911 gb AAQ56721.1 	serologically defined colon canc...		33	6.5
gi 24662899 ref NP_648506.1 	CG14131-PA [Drosophila melanog...		33	9.2

Alignments

>[gi|55238289|gb|EAA11136.3|](#) ENSANGP000000021866 [Anopheles gambiae str. PEST]
>[gi|58387327|ref|XP_315479.2|](#) ENSANGP000000021866 [Anopheles gambiae str. PEST]
Length = 289

Score = 229 bits (584), Expect = 7e-59
Identities = 132/260 (50%), Positives = 180/260 (69%), Gaps = 12/260 (4%)

Query: 3 SSRDKRSGSIYKFRDALIKKIEKKFEEERLLRLQRLVVLGQANNEPSEEAEPSSALMLT 62
S DK +G IY+FRD L K +FE E++L+RLQR+ ++E + L++T
Sbjct: 5 SQHDKDTGKIYRFRDKLESKFNAEFEEAEQKLIQLQRE-----QEAKERLANELLMT 55

Query: 63 CDEEDPFPADPSNATTADVVKVLRITLSVQIEQLSCKQDTFQHQLDEBTSRVTKLDRRLNEV 122
C++EDP +P +AT ADV+K+L+T+S +I QL KQDTFQ ++ E K RV KL+ ++N
Sbjct: 56 CEEEDPTINPDSATIADVVMKMLKTVSRKIGQLGKQDTFQKEVCEAKHRVLKLEAKMNTT 115

Query: 123 LFMAEAVKDELMFDDFQKATEQIVTPVPWFEPQAVSNTELKDLDRHLATDEAYSRLNLT 182
L M E VKDE++ + EQ P P FEF+AVSNE E+ DLDHRLATDEA+ RNLT
Sbjct: 116 LSMTEKVKDEMLARYW--TPQDDPPEPGEFFFAVSNEQEMNDLDRHLATDEAFCRNLTN 173

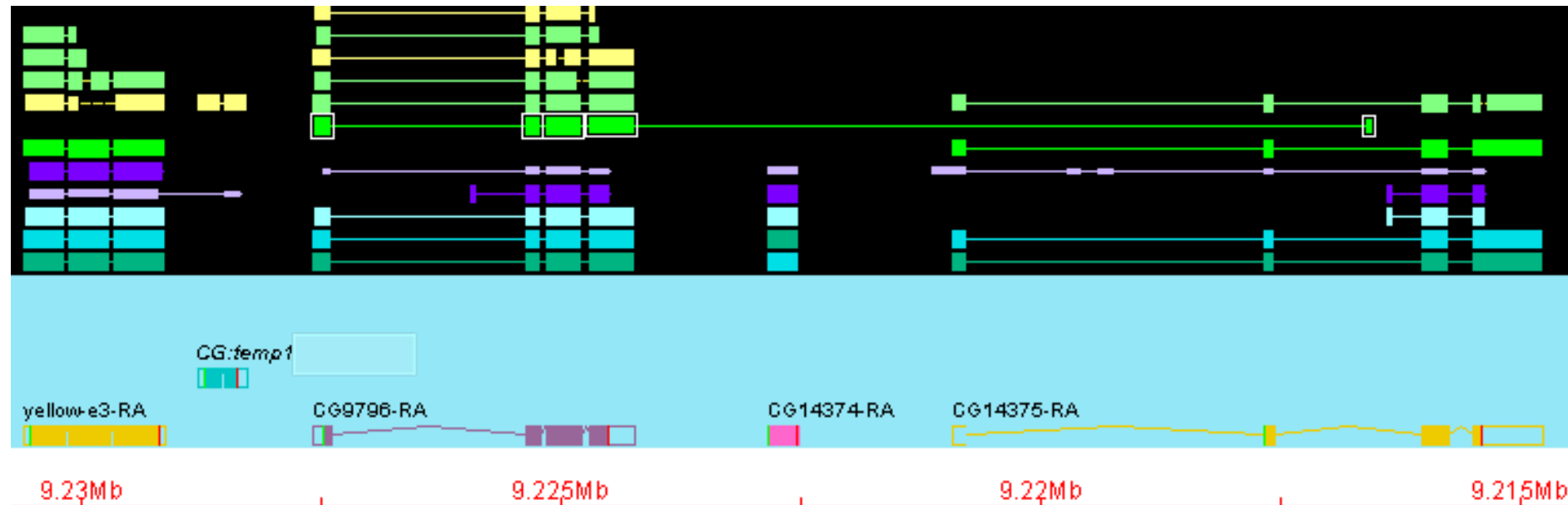
Query: 183 WLNAKILHPDPNHRHLQAMEAVFKREFLPSCSWKGRAPQGFRISMSAQKFIMELFRVVG 242
WL AKI P R+ A++ VF+REFL CSWGR+K G +I++ Q++I+ELFRVVG
Sbjct: 174 WLKAKISLQKPFRRMLHALDLVFEREFLMLCSWGRSKSGP-KIACDQRYIVELFRVVG 232

Query: 243 SNRFITISDEFVAKFFIKKL 262
SNR TI+D+ VAKFF+KLL
Sbjct: 233 SNRCTTITDKIVAKFFVRKL 252

CGtemp42 and CGtemp24 are
both homologues of ENSANG_21866



Manual annotation of Anopheles leads to annotation of new Drosophila genes



CLUSTAL W alignment
with DPM3 proteins

New Dmel gene
New Dpse gene
New Agam
gene

CLUSTAL W (1.82) multiple sequence alignment

```

DPM3-Bos      MTKLAQWLWALALLGSTWAALTMGALGLELPSSCREVLWPLPAYLLVSAG  50
DPM3-Mus      MTKLTQWLWGLALLGSAAWALTMGALGLELPFPCREVLWPLPAYLLVSAG  50
DPM3-Homo      MTKLAQWLWGLAILGSTWVALTTGALGLELPSCQEVWPLPAYLLVSAG  50
DPM3-Xenopus   MTKLAEWLLALSVLGAAWVTNLNGLGLDLPPLQQLWPLPVYLLVVF  50
CG_temp1-Dmel  MTNLQRWLFYASLFAIPYLSVVLGTQVPTLTKYFLHIQLPLLLLVIF  50
Dpse_Contig1107_Contig3589 MTNLQRWLFYATLFAVPYLSVVLGTQVPTLTKYFLHIQLPLLLLVIF  50
CG_temp4Anopheles MTKLFEWFMAAACFFSVYFAIVLRQVKHPLLDYMLBIQLSLPLFVLLF  50
                ***. *.: : : : : : : : : : * : * : *

DPM3-Bos      CYALGTVG YRVATPHDCEDAARELQSQIQEARADLTRRGLRF--  92
DPM3-Mus      CYALGTVG YRVATPHDCEDAARELQSQIVEARADLARRGLRF--  92
DPM3-Homo      CYALGTVG YRVATPHDCEDAARELQSQIQEARADLARRGLRF--  92
DPM3-Xenopus   CYSLATIG YRVATFNDCEAAARELQSQISEAKRDLALKGLKF--  92
CG_temp1-Dmel  IYSVWTVL YRTLTFTNDCPEAAKELQAEIQEARKDLIAKGFRFFD  94
Dpse_Contig1107_Contig3589 IYSVWTVL YRTTFTNDCPEAAKELQAEILEARKDLIAKGFRFFD  94
CG_temp4Anopheles IFSATVVL YRTTFTFNCEAAKELMEQIKEAKADLRSGKGLVLS  94
                :. :. *. *. :. :. :. :. :. :. :. :. :. :.
    
```

Results of *Anopheles gambiae* Manual Annotation



Initial Analysis of 1.9 Mb and 161 Ensembl annotations

- 49 Annotations required either additional coding exons, 5' or 3' sequence based on blast or EST or alternative transcripts
- 28 New genes annotated based on Gene prediction, ESTs and BLASTX
- 7 Gene merges - 14 Ensembl annotations merge into 7
- 10 Gene splits - 10 Ensembl annotations split into 22 new annotations
- 7 Ensembl annotations designated for deletion
- 25 Ensembl annotations required only minor adjustment to 5' and/or 3' sequence
- 56 Ensembl annotations required no change.

2L 15450000 - 15650000
2L 17750000 - 18000000
2L 30990000 - 31240000
2R 14450000 - 14650000
2R 32210000 - 32460000
2R 59310000 - 59510000
3L 10900000 - 11120000
3L 11650000 - 11900000

Results of *Anopheles gambiae* Manual Annotation



Initial Analysis of 1.9 Mb and 161 Ensembl annotations

- 50% of the genes predicted by the Ensembl Automatic Gene Building System required little or no change
- 30% of the genes required moderate changes to the gene structure
- 15% of automatic gene annotations were either split or merged (primarily based on BLAST evidence)
- The 28 new genes annotated comprised tandem related genes, nested genes, genes with low complexity sequence.

2L 15450000 - 15650000
2L 17750000 - 18000000
2L 30990000 - 31240000
2R 14450000 - 14650000
2R 32210000 - 32460000
2R 59310000 - 59510000
3L 10900000 - 11120000
3L 11650000 - 11900000

Limitations, challenges and the future of manual gene annotation



- While manual annotation provides the highest quality gene annotations, it is both labor intensive and expensive.
- The data illustrate the challenges for automated gene builds. Annotators will work together with Ensembl to rigorously analyze discrepancies between automated and manual Anopheles gene structures in an effort to improve gene prediction algorithms.
- These efforts will inform future automated annotation pipelines of not only Anopheles but also Aedes.
- Not unexpectedly, the initial pilot manual annotation Anopheles has led to corrections of several Drosophila melanogaster gene structures.

VectorBase



Ewan Birney
Martin Hammond
Vivek Iyer
Bill Gelbart
Frank Collins
Christos Louis
Fotis Kafatos

